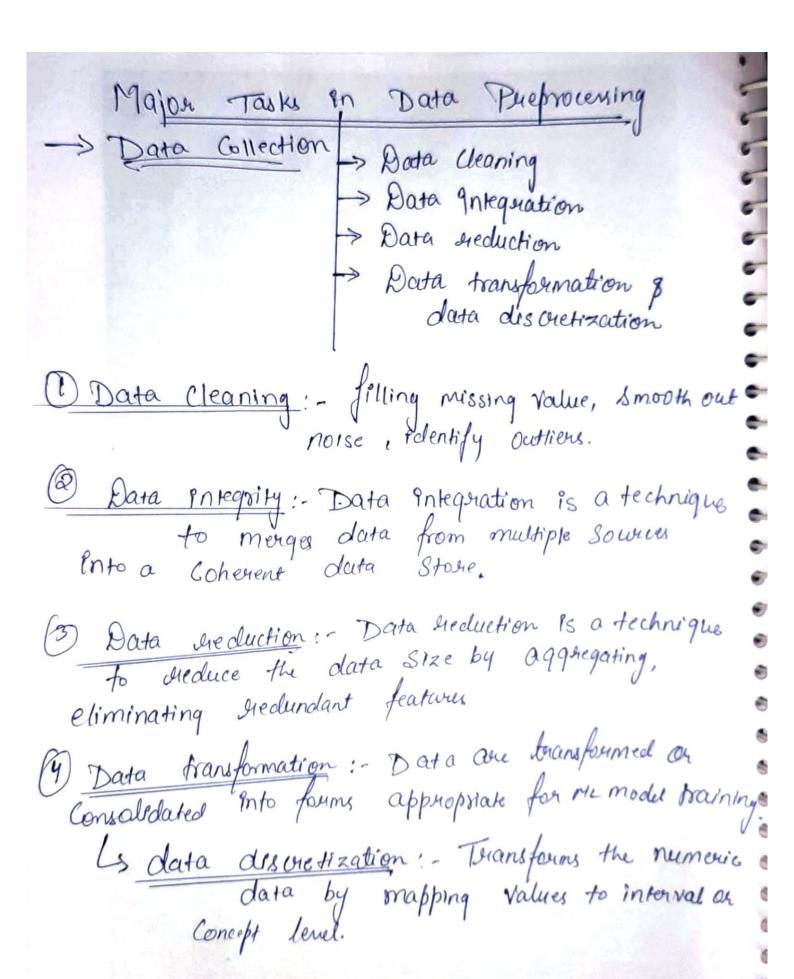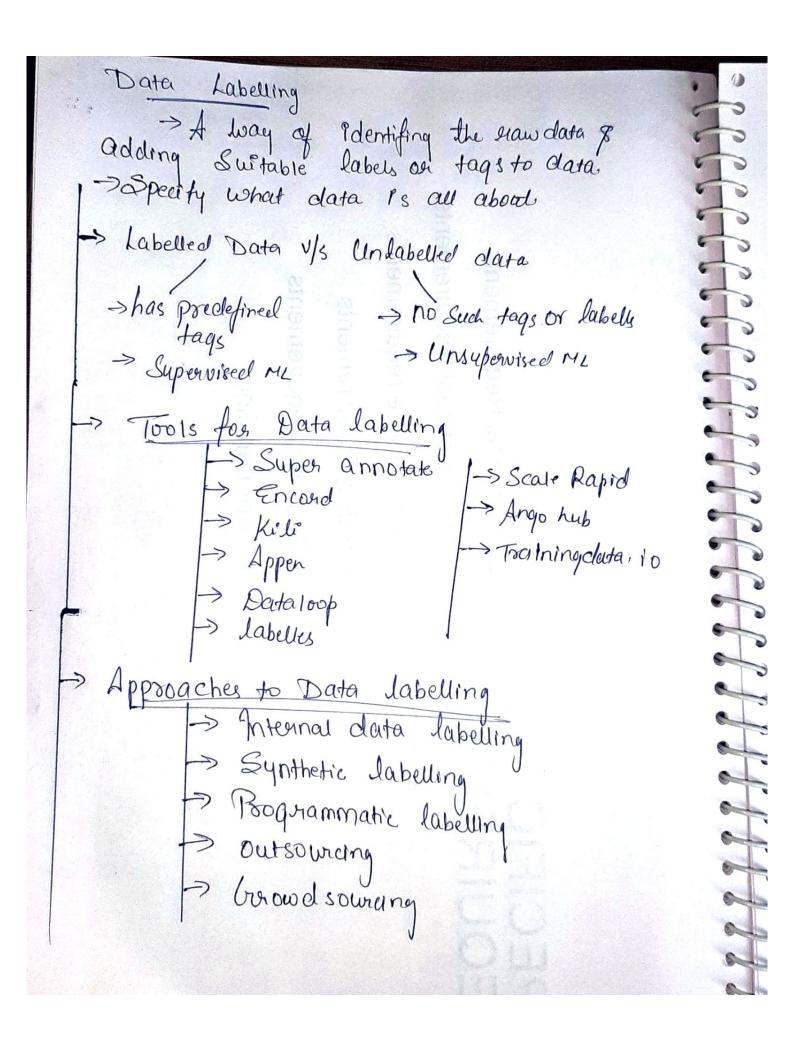# UNIT- 2 Data Pre-Processing

## Data Preprocessing : An Overview

→ Data Quality.

→ major tasks in Data Preprocessing

  → Data Cleaning
  → Data Integration
  → Data Reduction
  → Data Transformation & data discretization

Data pre-processing is the process of preparing the data & making it suitable for ML models.

→ Importance of data pre-processing
(Data in raw format contains noises, missing value, not in useable format)

→ Data Quality Measures

  → Accuracy
  → Completeness
  → Consistency
  → Timliness
  → Bellevability
  → Interpretability

# Major Tasks in Data Preprocessing

→ Data Collection
- → Data Cleaning
- → Data Integration
- → Data reduction
- → Data transformation & data discretization

① Data Cleaning :- filling missing value, Smooth out noise, identify outliers.

② Data Integrity :- Data Integration is a technique to merges data from multiple Sources into a Coherent data Store.

③ Data reduction :- Data reduction is a technique to reduce the data Size by aggregating, eliminating redundant features

④ Data transformation :- Data are transformed or Consolidated into forms appropriate for ML model training.

↳ data discretization :- Transforms the numeric data by mapping values to interval or concept level.

# Data Collection

Gathering data from various Sources Such as
- → Social media
- → databases
- → External repositories
- → IoT devices
- → Multimedia data

→ **Different data Collection methods**
- → Primary data Collection
- → Secondary data Collection

## Primary data Collection
- → Surveys & Questionnaries
- → Interviews
- → Observations
- → Experiments
- → focus Groups

## Secondary data Collection
- → Published Sources
- → Online Databases
- → Government & Institutional records
- → Publicly available data
- → Past research Studies

→ **Data Collection Tools**
- → Social media listening Tools
- → Web Analytics tools
- → Data Logging Devices
- → Mobile Data Coll'n apps
- → IoT devices

# Data Augmentation (A method of data coll$^n$)
(Expand the size of Existing dataset without gathering more data)

→ **Methods**

→ Rotating the Original image
⇒ Crop the original image differently
⇒ Altering the light conditions
→ Random Cropping and Padding
→ Scaling & Zooming
→ Shearing & Perspective transform
→ Colour Segmentation
→ Gaussian noise

→ **Types of Data Augmentation**

**Real data augmentation approches**

→ Sensor noise
→ Occlusion
→ Weather
→ Time Series
→ Label Smoothing

**Synthetic data augmentation approches**

→ Image Synthesis
→ Text generation
→ Oversampling and Undersampling
→ Data Interpolation & Extrapolation
→ feature Perturbation

→ **Challenges faced by Data augmentation**

→ Data Security & privacy
→ Maintaining label integrity
→ May increase size of training dataset

4

# Data Labelling

→ A way of identifying the raw data & adding suitable labels or tags to data.

→ Specify what data is all about

→ Labelled Data v/s Unlabelled data

→ has predefined tags

→ Supervised ML

→ no such tags or labells

→ Unsupervised ML

→ Tools for Data labelling

→ Super annotate
→ Encord
→ Kili
→ Appen
→ Dataloop
→ labelles

→ Scale Rapid
→ Ango hub
→ Trainingdata.io

→ Approaches to Data labelling

→ Internal data labelling
→ Synthetic labelling
→ Programmatic labelling
→ Outsourcing
→ Crowdsourcing

# Benifits of Data labelling

→ Precise Predictions

→ Better Data Usability

# Challanges of Data labelling

→ Costly & time Consuming

→ Possibilities of Human-errors

# Data Cleaning

(Data cleaning is the process of correcting or removing inacurate, impropely formated, duplicated or deal with missing values from dataset)

In general data cleaning lowers errors, improves the quality of the data

→ Steps for Cleaning data

→ Remove duplicate or Irrelevant observations
→ fix Structural errors
→ fitter unwanted Outliers
→ Handle missing data
→ Validate & QA

→ Techniques for Data cleaning

→ Ignore the tuples
→ fill in the missing value
→ Binning method
→ Regression
→ Clustering

→ **Process of Data Cleaning**

- → Monitoring the errors
- → Standardize the minning process
- → Validate data accuracy
- → Scrub for duplicate data
- → Research on data
- → Communicate with the team

→ **Tools for Data Cleaning**

- → OpenRefine
- → Trifacta Wrangler
- → Drake
- → Data Ladder
- → Cloudingo
- → Reifier (any many more)

→ **Benifits of Data cleaning**

- → Removes inaccuracies
- → Capacity to map many functions
- → Monitoring Mistakes
- → Makes decisions more quickly.
- → Enhance the Efficiency.

# Data Integration

(Process of merging data from several disparate sources)

→ Characteristics (G, S, M)

    G → global schema

    S → heterogeneous source of schema

    M → mapping b/w sources & global schema

→ Data Integration Approaches

    → Tight Coupling

    → low Coupling

→ Issues in Data Integration

    → Entity Identification Problem

    → Redundancy & Correlation Analysis

    → Tuple Duplication

    → Data warfare Detection & backbone

→ Data Integration Techniques

    → Manual Integration

    → Middleware Integration

    → Application- based Integration

    → Uniform access Integration

    → Data warehousing

→ Integration Tools

    → On- promise data Integration tool

    → open- source data integration tool

    → Cloud- based data integration tool

# Data Reduction

(Process of reducing the volume of Original data & represents it in a much Smaller volume)

→ By reducing the data, the efficiency of data analysis process is improved, which produces the same analytical results.

→ Importance of data reduction

  → Aims to define it in more compact form.
  → In terms of rows & colmu.

→ Techniques of Data Reduction

  → Dimensionality Reduction
  → Numerosity Reduction
  → Data Cube Aggregation
  → Data Compression

(I) → Dimensionality Reduction

(Eliminates attributes from data set by reducing vol. of original data)

  → A way of converting higher dimensions into lesser dimensions in a dataset.

  → Methods of Dimensionality Reduction

    → Feature Selection
    → Feature Extraction
    → Wavelet Transform
    → Attribute Subset Selection

# Dimensionality Reduction

**Feature Selection**
- → Filter method
- → Wrapper method
- → Embedded method

**Feature Extraction**
- → Principal Component Analysis (PCA)
- → Factor Analysis
- → Singular Value decomposition

→ **PCA** — Principal Component Analysis
- → Karl, Pearson, 1901
- → High dimensions → lower dimensions
- → Preserving most important patterns & relation between Variables

$n$ ⟶ Total no. of dimensions in a dataset

$k$ — independent features

$\underline{n > k}$

→ **Factor Analysis**
(There will not be any outliers, reduce a large no. of variables into fewer no. of factors)

→ **Singular value decomposition**
(helps us to simplify the data, finds most important patterns in the data & focus on them.)

# Wavelet Transformation

A data vector A is transformed into a numerically different data vector A' such that both A & A' vectors are of same length.

## Attribute Subset Selection

Reduces the volume of data by eliminating redundant & irrelevant attributes

## (II) Numerosity Reduction

→ Reduces the data volume by choosing alternative smaller forms of data representation

→ Types of numerosity reduction
  → Parametric numerosity reduction
  → Non-Parametric numerosity reduction

## (III) Data Cube Aggregation

→ This technique is used to aggregate data in a simpler form

→ It is a multidimensional aggregation that uses aggregation at various levels of a data cube to represent the original data set, thus achieving data reduction

## (iv) Data Compression

→ Employs modification, encoding & converting the structure of data in a way that consumes less space.

→ Types of Data Compression
  → Lossless data compression
  → Lossy data compression

# Data Transformation

→ Technique used to convert raw data into a suitable format that efficiently ease data analysis process.

→ **Importance**
  → Increases the efficiency
  → enables business to make better data-driven decisions.

→ **Data transformation may be:-**
  → Constructive
  → Destructive
  → Aesthetic
  → Structural

→ **Data Transformation Techniques**
  → Data Smoothing
  → Data Normalization
  → Data Discretization
  → Data Generalization

① **Data Smoothing**
  (Process used to remove noise from dataset)

  → Techniques to Remove Noise
    → Binning
    → Regression
    → Clustering

② Data Normalization
  (refers to scale the data values to a much smaller range)
  → Methods to Normalize the data
      → Min-max normalization
      → Z-score normalization
      → Decimal Scaling

③ Data Discretization
  (process of converting continuous data into a set of data intervals.)
  → Classification of data discretization
      → Supervised discretization
      → Unsupervised discretization

④ Data Generalization
  (converts low-level data attributes to high-level data attributes)
  → Divided into two approches
      → Data Cube process (OLAP)
      → Attribute-oriented Induction (AOI)

→ Data Transformation process
  (Entire process for transforming data is known as ETL)
      → Data discovery
      → Data Mapping
      → Data Extraction
      → Code Generation & Execution
      → Review
      → Sending.

# Advantages of Data Transformation

- → Better Organization
- → Improved data quality
- → Perform faster Queries
- → Better data Management
- → More Use out of Data

# Disadvantages of Data Transformation

- → Can be Expensive
- → Can be resource-Intensive
- → Lack of Expertise can Introduce problem

# Tools for Data Transformation

- → Scripting
- → On-premises ETL tools
- → Cloud-Based ETL tools